

Generalized Triangular Decomposition in Transform Coding

Ching-Chih Weng, *Student Member, IEEE*, Chun-Yang Chen, *Student Member, IEEE*, and P. P. Vaidyanathan, *Fellow, IEEE*

Abstract—A general family of optimal transform coders (TCs) is introduced here based on the generalized triangular decomposition (GTD) developed by Jiang *et al.* This family includes the Karhunen–Loeve transform (KLT) and the generalized version of the prediction-based lower triangular transform (PLT) introduced by Phoong and Lin as special cases. The coding gain of the entire family, with optimal bit allocation, is equal to that of the KLT and the PLT. Even though the original PLT introduced by Phoong *et al.* is not applicable for vectors that are not blocked versions of scalar wide sense stationary processes, the GTD-based family includes members that are natural extensions of the PLT, and therefore also enjoy the so-called MINLAB structure of the PLT, which has the unit noise-gain property. Other special cases of the GTD-TC are the geometric mean decomposition (GMD) and the bidiagonal decomposition (BID) transform coders. The GMD-TC in particular has the property that the optimum bit allocation is a uniform allocation; this is because all its transform domain coefficients have the same variance, implying thereby that the dynamic ranges of the coefficients to be quantized are identical.

Index Terms—Bit allocation, generalized triangular decomposition, geometric mean decomposition, linear prediction, majorization, Schur convexity.

I. INTRODUCTION

IN transform coder (TC) theory, the Karhunen–Loeve transform (KLT) is known for its optimality properties [2], [10], [25]. For example, it provides maximum coding gain when high-bit-rate scalar quantizers are used in the transform domain. The KLT essentially diagonalizes the autocorrelation matrix of the input vector \mathbf{x} before quantization. The decorrelated components are typically quantized by independent scalar quantizers.¹

If the vector \mathbf{x} being transformed is a blocked version of a scalar wide sense stationary (WSS) process $x(n)$, then the

coding gain of the KLT can also be achieved by using a different kind of transform called the prediction-based lower triangular transform (PLT), which was introduced into the signal-processing literature by Phoong and Lin [19]. The PLT is based on the theory of linear prediction for the scalar WSS process $x(n)$. PLT has smaller design cost because fast algorithms such as the Levinson algorithm can be used instead of matrix diagonalization. The implementation complexity for the PLT is 50% smaller than that of the [19]. However, the PLT as introduced in [19] is in the context of blocked versions of scalar WSS processes only, which is not applicable for general WSS vectors processes.

This paper introduces a general family for transform coding based on the generalized triangular decomposition (GTD) introduced by Jiang *et al.*, in the context of optimal transceiver design in digital communications [13]. We will show that the GTD-TC family has the following features.

- 1) Unlike the original PLT,² the input vector \mathbf{x} is not required to be a blocked version of a WSS process, but when such is the case the complexity of the new transform can be made comparable to that of the original PLT. One of the attractive features of the PLT is the existence of a structure with unit noise gain, called the MINLAB structure [19]. The GTD-based family includes a PLT-like special case which also enjoys the MINLAB structure. In this sense, it extends some of the features of the PLT for the case where \mathbf{x} is not a blocked version of a scalar process.
- 2) It includes the KLT and PLT as special cases.
- 3) The coding gain for any member of the family is equal to that of the KLT.
- 4) Like the KLT and the PLT, the GTD family also produces a decorrelated set of components at the inputs of the scalar quantizers. The GTD offers a great deal of freedom in the distribution of the variances of these decorrelated transform domain components.
- 5) Other special cases of the GTD transform coder includes the geometric mean decomposition (GMD) and the bidiagonal decomposition (BID) transform coders.
- 6) The GMD-TC in particular has the property that the optimum bit allocation is a uniform allocation. This follows from the fact that all transform coefficients have the same

Manuscript received March 16, 2009; accepted August 18, 2009. First published September 09, 2009; current version published January 13, 2010. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Marcelo G. S. Bruno. This work was supported in part by the Office of Naval Research under Grant N00014-08-1-0709 and the National Science Council, Republic of China, under Taiwan TMS scholarship 94-2-A-018.

The authors are with the California Institute of Technology, Pasadena, CA 91125 USA (e-mail: cweng@caltech.edu; cych@caltech.edu; ppvnath@systems.caltech.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSP.2009.2031733

¹Depending on the objective function to be optimized and the statistical assumptions involved, it can also be argued that the KLT is suboptimal in some sense [3]; we do not get into this aspect here.

²The original PLT, as introduced in [19], assumes that the input vector \mathbf{x} is a blocked version of a scalar WSS process. The natural extension of the PLT, introduced in Section II-B, will be shown to be optimal in terms of coding gain for any stationary vector process (not necessarily a blocked version of a scalar process) with well-defined covariance matrix [9]. This generalization will be referred to as the “PLT” and the restricted one in [19] as the “original PLT” throughout this paper.

variance (same dynamic range from a practical view point [23]), and thus the same machine word length can be used for all coefficients. Recall here that the closed-form formula for optimal bit allocation used by KLT and other transforms [10] often yields noninteger values for the bits. The approximation of these with integers would lead to suboptimality of the transform coder. Since the GMD-based method uses the same number of bits for all the transform domain coefficients without compromising optimality, this disadvantage is not present anymore.

The family of GTD coders therefore provides a unified framework for a number of optimal linear transforms for high-bit-rate coders.

This paper is organized as follows. Section II briefly reviews the KLT and the PLT. In Section III, we discuss the proposed GTD-TC. Several examples of the GTD-TC, such as the GMD-TC and BID-TC, are given here. The use of GTD in progressive transmission will also be described. Section IV provides numerical simulations related to the topic discussed in the paper. In particular, the theoretical claim that the GMD-TC with uniform bit allocation is as good as the KLT with optimal bit allocation is clearly demonstrated in this section. Section V concludes this paper.

Notations: Boldface uppercase letters denote matrices, boldface lower case letters denote column vectors, and italics denotes scalars. The superscripts $(\cdot)^T$ and $(\cdot)^\dagger$ denote the transpose and conjugate transpose operations. A_{ij} denotes the (i, j) th element of the matrix A . By $A \succeq B$, we mean that $A - B$ is positive semidefinite. For vector \mathbf{x} , the notation $\text{diag}(\mathbf{x})$ denotes the diagonal matrix with diagonal terms equal to the elements in the vector \mathbf{x} . For matrix \mathbf{X} , the notation $\text{diag}(\mathbf{X})$ denotes the column vector whose elements are the diagonal terms of \mathbf{X} . The notation $\mathbf{a} \prec_+ \mathbf{b}$ means that the vector \mathbf{b} majorizes \mathbf{a} additively [18], [16]. Similarly, $\mathbf{a} \prec_\times \mathbf{b}$ means that the vector \mathbf{b} majorizes \mathbf{a} multiplicatively [13], [16].

Assumptions: All signals and transforms discussed in this paper are assumed to be real-valued. We assume that the $M \times 1$ input $\mathbf{x}(n)$ is a zero-mean real-valued WSS vector process, with positive definite covariance matrix \mathbf{R}_x . The time argument n is dropped when redundant.

II. PRELIMINARIES AND REVIEWS

The transform coder is shown in Fig. 1. The signal \mathbf{x} is first multiplied by an $M \times M$ matrix \mathbf{T} so that $\mathbf{y} = [y_1 \ y_2 \ \cdots \ y_M]^T = \mathbf{T}\mathbf{x}$. The quantizers are scalar quantizers and are modeled as additive noise sources so that $\hat{y}_i = y_i + q_i$. Suppose the i th quantizer Q_i has b_i bits; then the variance of the quantization error q_i satisfies

$$\sigma_{q_i}^2 = c2^{-2b_i}\sigma_{y_i}^2 \quad (1)$$

where $\sigma_{y_i}^2$ is the variance of the signal input to the i th quantizer. This result generally holds under the high-bit-rate assumption [10], [17], [25]. The constant c depends on the type of quantizer and the statistics of y_i . It is assumed that all the scalar quantizers have the same c . The signal is reconstructed at the decoder by multiplying with \mathbf{T}^{-1} .

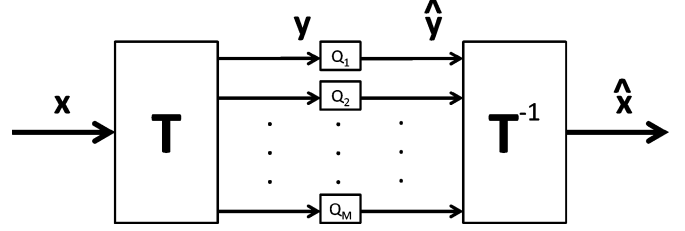


Fig. 1. Schematic of a transform coder with scalar quantizers.

A. Transform Coders and the KLT

The problem of minimizing the arithmetic mean of mean squared error (AM-MSE) of the reconstructed coefficients $E[\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2]$, under the average bit-rate constraint, is solved by the KLT [25]. The KLT uses $\mathbf{T} = \mathbf{U}^T$, where \mathbf{U} is any $M \times M$ orthonormal matrix such that $\mathbf{R}_x = \mathbf{U}\mathbf{\Sigma}\mathbf{U}^T$, and $\mathbf{\Sigma}$ is the diagonal matrix of the eigenvalues $\{\sigma_1^2, \dots, \sigma_M^2\}$ of \mathbf{R}_x (assumed to be in nonincreasing order).

Under the high-bit-rate assumption (1), the optimal bit allocation is given by the bit-loading formula [10], [25]

$$b_i = b + \frac{1}{2} \log_2 \frac{\sigma_i^2}{\det(\mathbf{R}_x)^{\frac{1}{M}}} \quad (2)$$

where the average bit rate is constrained to be b bits per data stream. Note that σ_i^2 is actually the signal variance of the transform coefficient y_i . With the bit allocation chosen as in (2), the MSE $\sigma_{q_i}^2$ due to the i th quantizer becomes independent of i (as seen by substituting (2) into (1), with $\sigma_{y_i} = \sigma_i$). The resulting AM-MSE is

$$\mathcal{E}_{\text{KLT}} = c2^{-2b} \det(\mathbf{R}_x)^{\frac{1}{M}}. \quad (3)$$

It was shown in [24] that under the high-bit-rate assumption, it is not a loss of generality to assume that the transform is orthonormal.³ It should be noted that the KLT decorrelates the signal, so the components of \mathbf{y} are statistically independent (under the Gaussian assumption) [4]. This is a necessary condition for optimality (minimum MSE) under the use of scalar quantizers [10] in the high-bit-rate case.

B. Prediction-Based Lower Triangular Transform (PLT)

The PLT, proposed in [19], is a signal dependent nonorthonormal transform, which utilizes linear prediction theory [10], [26]. It has the same decorrelation property as the KLT and is shown to have the same minimum MSE performance if the “minimum noise structure” and optimal bit allocation are used [19]. In [19], the original PLT is used for the vector \mathbf{x} obtained by blocking a scalar WSS $x(n)$. In the following review of the PLT idea, it can be seen that the PLT can actually be used for a vector process that need not be a blocked version of a scalar process. The development of [19], which was based on linear prediction theory, does not apply in

³It should be noted that the KLT is optimal among memoryless transforms. If \mathbf{T} is replaced with $\mathbf{T}(z)$, which has memory, then the lapped transform and its variations can be used to further improve the coding performance [1], [15]. In the lapped transform the optimal transform is no longer necessarily orthogonal but biorthogonal [15]. Such transforms are popular in modern practical transform coders [23].

- 3) the QR decomposition $\mathbf{H} = \mathbf{QR}$, where \mathbf{R} is an upper triangular matrix (here $\mathbf{P} = \mathbf{I}$);
- 4) the complete orthogonal decomposition $\mathbf{H} = \mathbf{Q}_2 \mathbf{R}_2 \mathbf{Q}_1^\dagger$, where $\mathbf{H}^\dagger = \mathbf{Q}_1 \mathbf{R}_1$ is the QR factorization of \mathbf{H}^\dagger and $\mathbf{R}_1^\dagger = \mathbf{Q}_2 \mathbf{R}_2$ is the QR factorization of \mathbf{R}_1^\dagger ;
- 5) the bidiagonal decomposition (BID) $\mathbf{H} = \mathbf{QRP}^\dagger$, where \mathbf{R} is a bidiagonal and upper triangular matrix [5, p. 251];
- 6) the geometric mean decomposition (GMD) [11], $\mathbf{H} = \mathbf{QRP}^\dagger$, where \mathbf{R} is an upper triangular matrix with the diagonal elements equal to the geometric mean of the positive singular values.

Now consider the transform coding problem again. Suppose the LDU decomposition of \mathbf{R}_x is $\mathbf{R}_x = \mathbf{LDL}^T$, as in (4). Decompose $\mathbf{D}^{\frac{1}{2}} \mathbf{L}^T$ using the GTD, i.e.,

$$\mathbf{D}^{\frac{1}{2}} \mathbf{L}^T = \mathbf{QRP}^T. \quad (9)$$

Then we can express \mathbf{R}_x as

$$\begin{aligned} \mathbf{R}_x &= \mathbf{PR}^T \mathbf{Q}^T \mathbf{QRP}^T \\ &= \mathbf{PL}_1 \text{diag}([\mathbf{R}_{11}^2, \mathbf{R}_{22}^2, \dots, \mathbf{R}_{MM}^2]) \mathbf{L}_1^T \mathbf{P}^T \end{aligned}$$

where \mathbf{L}_1 is a unit-diagonal lower triangular matrix that satisfies

$$\mathbf{L}_1 \text{diag}([\mathbf{R}_{11}, \mathbf{R}_{22}, \dots, \mathbf{R}_{MM}]) = \mathbf{R}^T.$$

Note that because of the GTD theory, the multiplicative majorization property

$$[\mathbf{R}_{11}^2, \mathbf{R}_{22}^2, \dots, \mathbf{R}_{MM}^2] \prec \times [\sigma_1^2, \sigma_2^2, \dots, \sigma_M^2] \quad (10)$$

holds, where $[\sigma_1^2, \sigma_2^2, \dots, \sigma_M^2]$ are the eigenvalues of \mathbf{R}_x with nonincreasing order, i.e., $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_M^2$. Note that (10) implies the fact that the diagonal terms of \mathbf{R} cannot be arbitrarily chosen but have to satisfy the multiplicative majorization property.

If we pass the signal \mathbf{x} through the orthonormal matrix \mathbf{P}^T to produce \mathbf{z} , i.e., $\mathbf{z} = \mathbf{P}^T \mathbf{x}$, the covariance of \mathbf{z} is

$$\mathbf{R}_z = \mathbf{P}^T \mathbf{R}_x \mathbf{P} = \mathbf{L}_1 \text{diag}([\mathbf{R}_{11}^2, \mathbf{R}_{22}^2, \dots, \mathbf{R}_{MM}^2]) \mathbf{L}_1^T.$$

Therefore, \mathbf{L}_1 is the lower triangular matrix of the LDU form of \mathbf{R}_z . If now apply the PLT \mathbf{L}_1^{-1} to the signal \mathbf{z} , the components of the resulting vector are decorrelated. The system is called GTD-TC and is demonstrated in Fig. 4 for $M = 4$. Here we have used the MINLAB(I) structure [19]. The multipliers s_{km} are the entries of the matrix \mathbf{L}_1^{-1} . For example, when $M = 4$

$$\mathbf{L}_1^{-1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ s_{21} & 1 & 0 & 0 \\ s_{31} & s_{32} & 1 & 0 \\ s_{41} & s_{42} & s_{43} & 1 \end{bmatrix}.$$

The bit-loading formula becomes

$$\begin{aligned} b_i &= b + \frac{1}{2} \log_2 \frac{\mathbf{R}_{ii}^2}{\det(\mathbf{R}_z)^{\frac{1}{M}}} \\ &= b + \frac{1}{2} \log_2 \frac{\mathbf{R}_{ii}^2}{\det(\mathbf{R}_x)^{\frac{1}{M}}} \end{aligned} \quad (11)$$

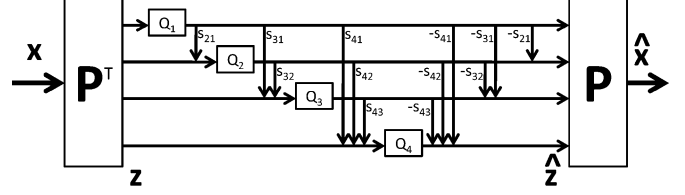


TABLE I
DESIGN AND IMPLEMENTATION COSTS OF TRANSFORM CODERS

	Design cost	Impl. cost(precoder part)	Impl. cost(PLT part)
KLT	EVD, $O(M^3)$	$O(M^2)$	0
PLT	LDU, $O(M^2)$	0	$O(M^2)$
GMD-TC	EVD and GMD [13], $O(M^3)$	$O(M^2)$	$O(M^2)$
BID-TC	Hessenberg form $O(M^3)$ and easy LDU $O(M)$	$O(M^2)$	$O(M)$
General GTD-TC	EVD and GTD [13], $O(M^3)$	$O(M^2)$	$O(M^2)$

because $\det(\mathbf{R}_x) = \bar{\sigma}^{2M}$. The preceding equation says that all the quantizers are assigned the same number of bits. This is a consequence of the fact that D_{ii} in (5) are identical for all i . That is, the variances of the quantizer inputs are all identical, which means that the dynamic ranges of the signals being quantized are identical. This is a desirable property in practice.

B. Bidiagonal Transformation—Hessenberg Form

A matrix \mathbf{B} is said to be bidiagonal if it has the form demonstrated below for the 4×4 case

$$\mathbf{B} = \begin{bmatrix} b_{00} & b_{01} & 0 & 0 \\ 0 & b_{11} & b_{12} & 0 \\ 0 & 0 & b_{22} & b_{23} \\ 0 & 0 & 0 & b_{33} \end{bmatrix}.$$

If the GTD form of $\mathbf{D}^{1/2}\mathbf{L}^T$ is \mathbf{QBP}^T , where \mathbf{B} is a bidiagonal matrix, then we call it the bidiagonal transform coder (BID-TC). It can be seen that

$$\mathbf{R}_x = \mathbf{LDL}^T = \mathbf{PB}^T\mathbf{BP}^T$$

where $\mathbf{B}^T\mathbf{B}$ is a tridiagonal matrix demonstrated below for size 4×4

$$\mathbf{B}^T\mathbf{B} = \begin{bmatrix} c_{00} & c_{01} & 0 & 0 \\ c_{10} & c_{11} & c_{12} & 0 \\ 0 & c_{21} & c_{22} & c_{23} \\ 0 & 0 & c_{32} & c_{33} \end{bmatrix}$$

with $c_{mk} = c_{km}$. This tridiagonal form $\mathbf{B}^T\mathbf{B}$ is also known as the Hessenberg form [5] of \mathbf{R}_x .

The advantages of the BID-TC coder lie in its reduced computational complexity. To reduce a symmetric matrix to a tridiagonal form by orthonormal transformation is computationally much less complex compared to eigenvalue decomposition [5]. The detail of reducing a symmetric matrix to the tridiagonal form is discussed in [5] and requires only several Householder transformations. The LDU decomposition for a symmetric tridiagonal matrix is also easy, and requires only $O(M)$ operations now instead of $O(M^2)$ for general symmetric matrices. Therefore, the design cost for the BID-TC is less than KLT, whereas KLT requires iterative EVD computations. Also, due to the bidiagonal structure of \mathbf{B} , the implementation cost for the inner PLT part is also reduced, which is only on the order of $O(M)$. This can be seen in Fig. 5, which shows the MINLAB(I) structure for the BID-TC encoder. Signal feedforward paths are only required

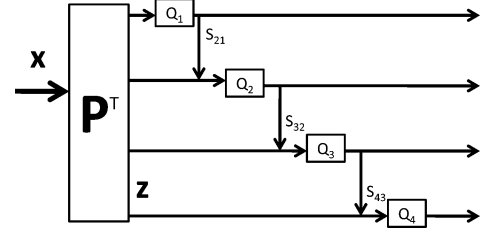


Fig. 5. The BID transform coder implemented using MINLAB(I) structure.

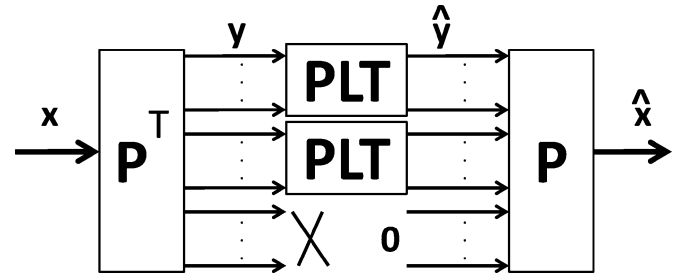


Fig. 6. Use of GTD-TC in the progressive transmission context.

for the adjacent data streams. The number of signal feedforward paths is much less than for the original PLT.

A detailed comparison between the design and implementation costs for various GTD-based coders is summarized in Table I.⁵

C. Combination of GMD and Progressive Transmission

There are some applications where rapid transmission is required and a coarse signal approximation is first produced [14]. When more bits are available, the system progressively enhances the performance by sending more information. Fig. 6 shows the example in which we divide the signal data streams after the linear transformation into three groups. The first group is the significant group, where the K_1 data streams contain a coarse approximation of the signal. The second group is the less significant group, where the K_2 data streams contain detailed information about the signal. The third group of K_3 streams is the least significant group, where the remaining $M - K_1 - K_2$ data streams contain components that are close to zero after the linear transformation \mathbf{P}^T .

⁵In situations involving the KLT, the discrete cosine transform (DCT) is often used instead of the KLT since the DCT is signal independent, computationally efficient, and a good approximation of the KLT for a large class of signals with low-pass spectra [15]. An analogous low-complexity approximation for the precoder \mathbf{P} that arises in the GTD implementation is not known at this time.

Suppose we adopt the GTD form in (9). We are looking for a transformation such that the diagonal terms of \mathbf{R} have the pattern

$$\text{diag}(\mathbf{R}) = [\bar{\sigma}_1, \dots, \bar{\sigma}_1, \bar{\sigma}_2, \dots, \bar{\sigma}_2, \sigma_{K_1+K_2+1}, \dots, \sigma_M]$$

where

$$\bar{\sigma}_1^2 = \left(\prod_{i=1}^{K_1} \sigma_i^2 \right)^{\frac{1}{K_1}}, \quad \bar{\sigma}_2^2 = \left(\prod_{i=K_1+1}^{K_1+K_2} \sigma_i^2 \right)^{\frac{1}{K_2}}.$$

Here $[\sigma_1^2, \dots, \sigma_M^2]$ are the eigenvalues of \mathbf{R}_x with nonincreasing order. \mathbf{P} is the orthonormal matrix obtained from the GTD theory. Note that this decomposition exists for any K_1, K_2 combination, since the multiplicative majorization property holds. Because the eigenvalues are in nonincreasing order, the first K_1 substreams actually represent the first K_1 principal components of the vector \mathbf{x} , and the next K_2 substreams represent the next K_2 principal components. Suppose for the significant group the total bit budget is $b_1 K_1$, for the less significant group the total bit budget is $b_2 K_2$, and for the least significant group the average number of bits are zero. As shown in Fig. 6, for the first and the second group, we use the local PLT for each of them. It can be seen that the bit-loading formula under the high-bit-rate assumption will be

$$b_i = b_1 + \frac{1}{2} \log_2 \frac{\mathbf{R}_{ii}^2}{\left(\prod_{i=1}^{K_1} \sigma_i \right)^{\frac{1}{K_1}}} = b_1$$

for the first group and

$$b_i = b_2 + \frac{1}{2} \log_2 \frac{\mathbf{R}_{ii}^2}{\left(\prod_{i=K_1+1}^{K_1+K_2} \sigma_i \right)^{\frac{1}{K_2}}} = b_2$$

for the second group. That is, uniform bit loading is used across the quantizers within each group. The data streams in the third group are dropped (i.e., assigned zero bits). It can be seen that the resulting AM-MSE of this transform coder is

$$\frac{1}{M} (K_1 c 2^{-2b_1} \bar{\sigma}_1 + K_2 c 2^{-2b_2} \bar{\sigma}_2 + \sum_{i=K_1+K_2+1}^M \sigma_i).$$

When we only have very low bit budget, we can allocate the bits to the first group to get the coarse approximation of the signal. When we have more bits available, the information in the second group is exploited to get the detailed information of the signal. Hence the progressive transmission scheme can be implemented when we are able to use uniform quantizers within each group. This shows one example of the flexibility that our proposed GTD-TC scheme can have. One can have more groups of data streams where each group has a different bit budget.

D. An Illustrative Example

Before we proceed to simulations, we provide a simple numerical example of the GMD and the BID transform coders for increased clarity of exposition. Suppose the zero-mean input vector \mathbf{x} is of size 3×1 , with covariance matrix

$$\mathbf{R}_x = \begin{pmatrix} 5.8607 & -1.2345 & -1.0333 \\ -1.2345 & 3.3396 & 0.3979 \\ -1.0333 & 0.3979 & 1.4943 \end{pmatrix}.$$

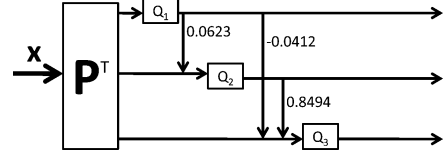


Fig. 7. The GMD encoder structure for the example.

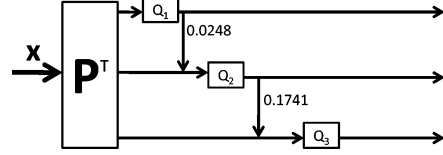


Fig. 8. The bidiagonal (BID) encoder structure for the example.

Following the procedures in Section III-A, the \mathbf{P} matrix can be calculated by the GMD, which is

$$\mathbf{P} = \begin{pmatrix} 0.2947 & 0.3478 & -0.8901 \\ 0.9556 & -0.1022 & 0.2765 \\ 0.0052 & -0.9320 & -0.3625 \end{pmatrix}.$$

The resulting \mathbf{L}_1 is a unit diagonal lower triangular matrix such that

$$\mathbf{L}_1^{-1} = \begin{pmatrix} 1.0000 & 0 & 0 \\ 0.0623 & 1.0000 & 0 \\ -0.0412 & 0.8494 & 1.0000 \end{pmatrix}.$$

The coefficients of \mathbf{L}_1^{-1} are the multipliers s_{km} in MINLAB structure implementation of the encoder structure, as shown in Fig. 7. The signals at the input to the quantizers become uncorrelated, with identical variances [2.8640, 2.8640, 2.8640]. This allows the optimal bit-allocation scheme to be uniform.

For the same signal \mathbf{x} , if we use the BID-TC, the orthonormal matrix \mathbf{P} can be shown to be

$$\mathbf{P} = \begin{pmatrix} 0.3593 & 0.9332 & 0 \\ 0.9332 & -0.3593 & 0 \\ 0 & 0 & 1.0000 \end{pmatrix}.$$

The resulting \mathbf{L}_1 is a bidiagonal matrix such that

$$\mathbf{L}_1^{-1} = \begin{pmatrix} 1.0000 & 0 & 0 \\ 0.0248 & 1.0000 & 0 \\ 0 & 0.1741 & 1.0000 \end{pmatrix}.$$

The bidiagonal structure of \mathbf{L}_1 makes the implementation of the transform coder less complex (Fig. 8). The signal components at the input to the quantizers become uncorrelated, with the *unequal* variances [2.8372, 6.3614, 1.3016]. The optimal bit loading is then calculated according to (5), which in general does not yield integer values. Replacing these with integers reduces the coding gain from its theoretical value.

IV. SIMULATIONS

In this section, we provide the numerical simulations for GTD-based coders. The signal \mathbf{x} is generated by a zero-mean Gaussian vector process with prescribed covariance matrix

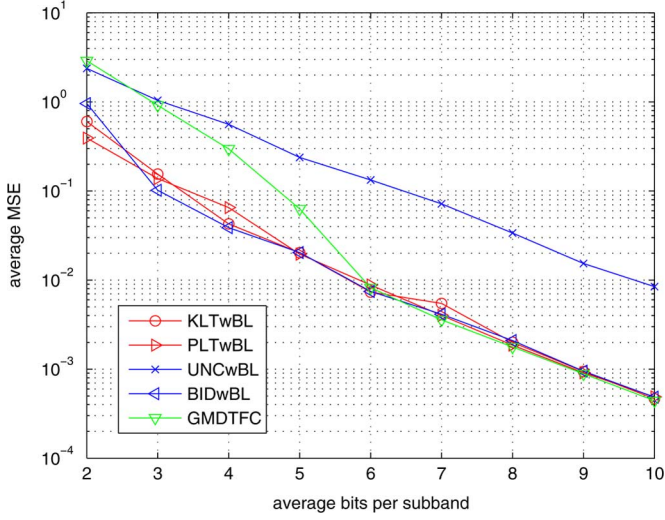


Fig. 9. Performance of different transform coders with optimal bit allocation. Input covariance matrix has high condition number (10^7).

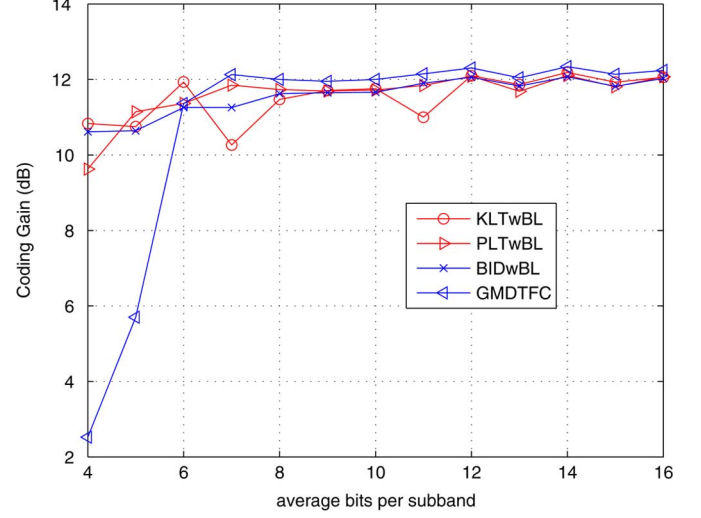


Fig. 11. Comparison of coding gain of different transform coders with optimal bit allocation. Input covariance matrix has high condition number (10^7).

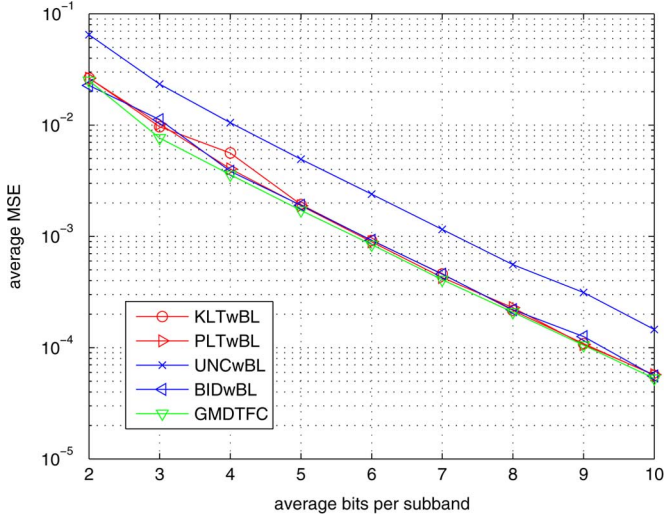


Fig. 10. Performance of different transform coders with optimal bit allocation. Input covariance matrix has low condition number (10^3).

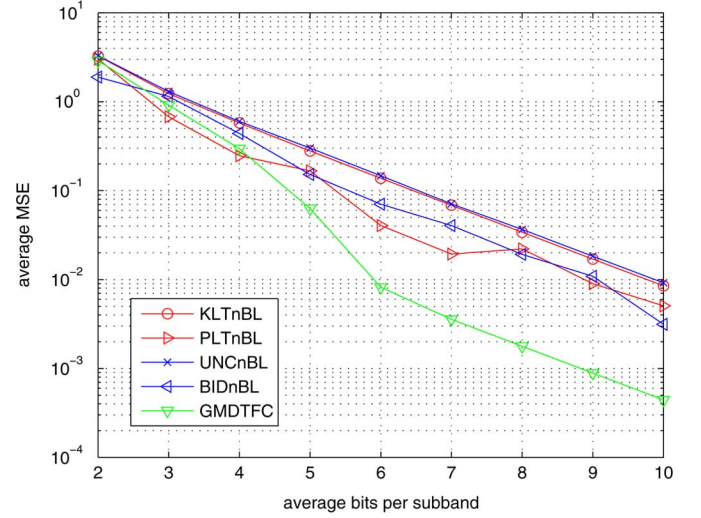


Fig. 12. Performance of different transform coders with uniform bit allocation. Input covariance matrix has high condition number (10^7).

\mathbf{R}_x . The number of data streams $M = 8$ in the experiments. Uniform roundoff quantizers are assumed. Each quantizer adapts its step size according to the variance of the Gaussian input [25, p. 818]. We run the simulation for input covariance with high and low condition numbers, respectively. In Figs. 9 and 12, the condition number is 10^7 . In Figs. 10 and 13, the condition number is 10^3 . For each case, we run the Monte Carlo simulations for calculating the AM-MSE. In each trial, we first generate the input covariance matrix by multiplying a fixed diagonal matrix $\mathbf{\Lambda}_x$ with a randomly generated orthonormal matrix on the left and its transpose on the right. Two choices of $\mathbf{\Lambda}_x$ are used. For the so-called high condition number example

$$\mathbf{\Lambda}_x = \text{diag}[10^7 \ 10^6 \ 10^5 \ 10^4 \ 10^3 \ 10^2 \ 10^1 \ 1]$$

and for the low condition number example

$$\mathbf{\Lambda}_x = \text{diag}[10^3 \ 10^3 \ 10^2 \ 10^2 \ 10^1 \ 10^1 \ 1 \ 1].$$

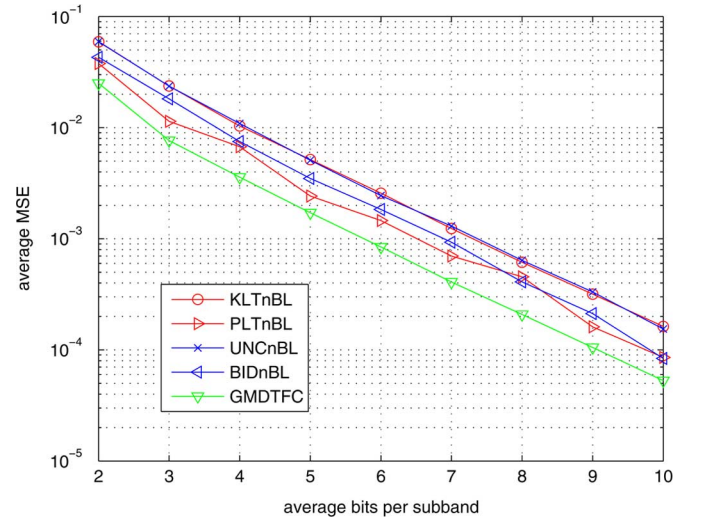


Fig. 13. Performance of different transform coders with uniform bit allocation. Input covariance matrix has low condition number (10^3).

The input vector \mathbf{x} is then generated according to this covariance matrix. In the following, we provide simulation comparisons of different transform coders with and without optimal bit allocation.

Optimal Bit Allocation: Figs. 9 and 10 compare the AM-MSE performance of different transform coders with optimal bit allocation for input covariance matrix with high and low condition numbers, respectively. “Transform-wBL” means we adopt the specified transform with optimal bit loading. For example, “KLTwBL” uses the KLT with the bit-loading formula (2). “PLTwBL” is the method mentioned in [19], with the optimal bit-loading formula (5). “UNCwBL” is the case when we have no transform; we directly quantize the input \mathbf{x} with optimal bit allocation⁶

$$b_i = b + \frac{1}{2} \log_2 \frac{\sigma_{x,i}^2}{(\prod_{k=1}^M \sigma_{x,k}^2)^{\frac{1}{M}}}.$$

Since the inputs to the quantizers x_i are correlated to each other in general, direct scalar quantization without transformation results in performance loss compared to the GTD-TCs even when the optimal bit loading scheme is applied. “BIDwBL” is the bidiagonal transform coder discussed in Section III-B. The bit-loading formula is as in (11). “GMDTFC” is the GMD transform coder. Since the signal variance in each data stream is the same, no bit loading is needed. This allows us to build the same scalar quantizers for all data streams. It can be seen from the figure that with optimal bit loading, all GTD-TCs perform about the same. This is consistent with the analysis made in Section III. Direct quantization without transforms (UNCwBL) results in about 5 bits per data stream performance loss for Fig. 9 and 1.7 bits loss for Fig. 10.

Fig. 11 plots the coding gain defined as

$$G_{TC} = \frac{\text{MSE}_{\text{PCM}}}{\text{MSE}_{\text{TC}}} \quad (14)$$

which is the ratio of the MSE of direct quantization MSE_{PCM} (often referred to as pulse coded modulation) to the MSE of the transform coder MSE_{TC} . It can be seen that the coding gain performance of each method is approximately the same in the high-bit-rate regime.

Uniform Bit Allocation: Figs. 12 and 13 compare the AM-MSE performance of different transform coders with uniform bit allocation for input covariance matrix with high and low condition numbers, respectively. Here, “transform-nBL” means we adopt some specific transform with no optimal bit loading, i.e., we allocate the same number of bits to each data stream. However, the step size of each scalar quantizer is adapted according to variance of the Gaussian input [25, p. 818]. “KLTnBL” uses KLT for the transform. “PLTnBL” is the method mentioned in [19] but with no bit loading. “UNCnBL” is the case when we have no transform but directly quantize the input \mathbf{x} . No bit loading is applied either. “BIDwBL” is the bidiagonal transform coder discussed in Section III-B with

⁶We perform a rounding operation on the bit-loading formula to obtain integer values and adjust it a little bit to fit the bit budget: first we check if the bit budget is satisfied with equality. If the number of bits is more/less than the bit budget, we decrease/increase 1 bit from the substream with most/least number of bits. We repeat this until the bit budget is satisfied with equality. While suboptimal, we believe this algorithm is not far from optimal in the high-bit-rate case.

no bit loading. “GMDTFC” is the GMD transform coder. It can be seen from the figure that with no bit loading applied, GMD performs much better than the other methods, since the GMD without bit allocation is theoretically as good as the other methods with optimal bit allocation.

In the simulation results, the reader will notice that for values of b (average number of bits) exceeding three (low condition number case) and exceeding six (for high condition number case), the theoretical predictions are indeed verified to be true.⁷ Namely, with no bit allocation, GMD performs much better than KLT, PLT, and BID. These latter methods with no bit allocation have performance comparable to direct quantization. Furthermore, with optimal bit allocation, all these methods (GMD, KLT, and BID) have identical performances. For small values of b [14], these theoretical predictions (which are based on the high-bit-rate assumption) are seen to be (understandably) less and less true. The low-bit-rate effect appears to be more severe for the case where the input covariance matrix has high condition number. Also, from the simulations, we see that the coding gain improvement of the proposed GTD-TC is more significant for the high condition number case.

V. CONCLUSION

The main purpose of this paper has been to provide a general framework for a family of linear transform coders based on the GTD. The GTD has in the past been found to be of great importance in digital transceiver optimization but has hitherto not been considered for transform coding. The KLT and PLT transforms are special cases belonging to the GTD transform coder family. Some of the new transform coders that have been presented as members of this family include the GMD and the BID coders. The BID has the advantage that the computational complexity of the PLT part is significantly less. The GMD has the special property that optimal bit allocation is actually a uniform allocation because all the transformed coefficients have identical variances or dynamic ranges.

REFERENCES

- [1] S. O. Aase and T. A. Ramstad, “On the optimality of nonunitary filter banks in subband coders,” *IEEE Trans. Image Process.*, vol. 4, pp. 1585–1591, Dec. 1995.
- [2] A. N. Akansu and Y. Liu, “On signal decomposition techniques,” *Opt. Eng.*, vol. 30, pp. 912–920, Jul. 1991.
- [3] M. Effros, H. Feng, and K. Zeger, “Suboptimality of Karhunen-Loeve transformation for transform coding,” *IEEE Trans. Inf. Theory*, vol. 50, pp. 1605–1619, Aug. 2004.
- [4] A. Gersho and R. M. Gray, *Vector Quantization and Signal Compression*. Norwell, MA: Kluwer Academic, 1992.
- [5] G. H. Golub and C. F. Van Loan, *Matrix Computations*. Baltimore, MD: Johns Hopkins Univ. Press, 1996.
- [6] R. A. Horn, “On the eigenvalues of a matrix with prescribed singular values,” *Proc. Amer. Math. Soc.*, vol. 5, no. 1, pp. 4–7, Feb. 1954.
- [7] R. A. Horn and C. R. Johnson, *Topics in Matrix Analysis*. Cambridge, U.K.: Cambridge Univ. Press, 1991.
- [8] R. A. Horn and C. R. Johnson, *Matrix Analysis*. Cambridge, U.K.: Cambridge Univ. Press, 1985.
- [9] J. J. Y. Huang and P. M. Schultheiss, “Block quantization of correlated Gaussian random variables,” *IEEE Trans. Commun. Syst.*, vol. CS-11, pp. 289–296, Sep. 1963.
- [10] N. S. Jayant and P. Noll, *Digital Coding of Waveforms*. Englewood Cliffs, NJ: Prentice-Hall, 1984.

⁷It should be mentioned here that such relatively large values for b are not uncommon in areas such as multispectral image compression [21].

- [11] Y. Jiang, J. Li, and W. W. Hager, "Joint transceiver design for MIMO communications using geometric mean decomposition," *IEEE Trans. Signal Process.*, pp. 3791–3803, Oct. 2005.
- [12] Y. Jiang, J. Li, and W. W. Hager, "Uniform channel decomposition for MIMO communications," *IEEE Trans. Signal Process.*, pp. 4283–4294, Nov. 2005.
- [13] Y. Jiang, W. W. Hager, and J. Li, "Generalized triangular decomposition," *Math. Comput.*, Oct. 2007.
- [14] S. Mallat and F. Falzon, "Analysis of low bit rate image transform coding," *IEEE Trans. Signal Process.*, vol. 46, pp. 1027–1042, Apr. 1998.
- [15] H. S. Malvar, "Biorthogonal and nonuniform lapped transforms for transform coding with reduced blocking and ringing artifacts," *IEEE Trans. Signal Process.*, vol. 46, pp. 1043–1053, Apr. 1998.
- [16] A. W. Marshall and I. Olkin, *Inequalities: Theory of Majorization and Its Applications*. New York: Academic, 1979.
- [17] D. L. Mary and D. T. M. Slock, "A theoretical high-rate analysis of causal versus unitary online transform coding," *IEEE Trans. Signal Process.*, vol. 54, pp. 1472–1482, Apr. 2006.
- [18] D. P. Palomar and Y. Jiang, "MIMO transceiver design via majorization theory," *Found. Trends Commun. Inf. Theory*, vol. 3, no. 4, pp. 331–551, Nov. 2006.
- [19] S. M. Phoong and Y. P. Lin, "Prediction-based lower triangular transform," *IEEE Trans. Signal Process.*, vol. 48, pp. 1947–1955, Jul. 2000.
- [20] S. M. Phoong and Y. P. Lin, "MINLAB: Minimum noise structure for ladder-based biorthogonal filter banks," *IEEE Trans. Signal Process.*, vol. 48, pp. 465–476, Feb. 2000.
- [21] J. A. Saghi, A. G. Tescher, and J. T. Reagan, "Practical transform coding of multispectral imagery," *IEEE Signal Process. Mag.*, pp. 32–43, Jan. 1995.
- [22] M. B. Shenouda and T. N. Davidson, "A framework for designing MIMO systems with decision feedback equalization or Tomlinson-Harashima precoding," *IEEE J. Sel. Areas Commun.*, vol. 26, pp. 401–411, Feb. 2008.
- [23] S. Srinivasan, C. Tu, S. L. Regunathan, R. A. Rossi Jr., and G. J. Sullivan, "HD photo: A new image coding technology for digital photography," in *Proc. SPIE Appl. Digital Image Process. XXX*, San Diego, CA, Aug. 2007, vol. 6696.
- [24] P. P. Vaidyanathan, "Theory of optimal orthonormal subband coders," *IEEE Trans. Signal Process.*, vol. 46, pp. 1528–1543, Jun. 1998.
- [25] P. P. Vaidyanathan, *Multirate Systems and Filter Banks*. Englewood Cliffs, NJ: Prentice-Hall, 1993.
- [26] P. P. Vaidyanathan, *The Theory of Linear Prediction*. London, U.K.: Morgan & Claypool, 2008.
- [27] C. C. Weng, C. Y. Chen, and P. P. Vaidyanathan, "MIMO transceivers with decision feedback and bit loading: Theory and optimization," *IEEE Trans. Signal Process.*, accepted for publication.
- [28] H. Weyl, "Inequalities between the two kinds of eigenvalues of linear transformations," in *Proc. Nat. Acad. Sci.*, July 1949, vol. 35, pp. 408–411.
- [29] J. K. Zhang, A. Kavcic, and K. M. Wong, "Equal-diagonal QR decomposition and its application to precoder design for successive-cancellation detection," *IEEE Trans. Inf. Theory*, pp. 154–172, Jan. 2005.



Ching-Chih Weng (S'07) was born in Taipei, Taiwan, R.O.C., on October 15, 1981. He received the B.S. degree from National Taiwan University, Taipei, in 2004 and the M.S. degree from the California Institute of Technology, Pasadena, in 2007, both in electrical engineering. He is currently pursuing the Ph.D. degree in digital signal processing at Caltech.

In summer 2009, he was with Qualcomm Inc. as an Intern System Engineer. His research interests include digital communications and radar signal processing.



Chun-Yang Chen (S'06) was born in Taipei, Taiwan, R.O.C., on November 22, 1977. He received the B.S. and M.S. degrees in electrical engineering and communication engineering from National Taiwan University (NTU), Taipei, in 2000 and 2002, respectively, and the Ph.D. degree in electrical engineering from the California Institute of Technology, Pasadena, in 2009.

His doctoral work was in the field of digital signal processing. He is currently with Facebook.



P. P. Vaidyanathan (S'80–M'83–SM'88–F'91) was born in Calcutta, India, on October 16, 1954. He received the B.Sc. (Hons.) degree in physics and the B.Tech. and M.Tech. degrees in radiophysics and electronics, all from the University of Calcutta, India, in 1974, 1977, and 1979, respectively, and the Ph.D. degree in electrical and computer engineering from the University of California at Santa Barbara in 1982.

He was a postdoctoral fellow at the University of California, Santa Barbara from September 1982 to March 1983. In March 1983, he joined the Electrical Engineering Department of the California Institute of Technology as an Assistant Professor, where he has been Professor of electrical engineering since 1993. His main research interests are in digital signal processing, multirate systems, wavelet transforms and signal processing for digital communications.

Dr. Vaidyanathan served as Vice-Chairman of the Technical Program committee for the 1983 IEEE International Symposium on Circuits and Systems, and as the Technical Program Chairman for the 1992 IEEE International Symposium on Circuits and Systems. He was an Associate Editor for the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS for the period 1985–1987, and is currently an Associate Editor for the IEEE SIGNAL PROCESSING LETTERS, and a consulting editor for the journal *Applied and Computational Harmonic Analysis*. He was a guest editor in 1998 for special issues of the IEEE TRANSACTIONS ON SIGNAL PROCESSING and the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS II, on the topics of filter banks, wavelets and subband coders. He has authored more than 360 papers in journals and conferences, and is the author of the book *Multirate Systems and Filter Banks*. He has written several chapters for various signal processing handbooks. He was a recipient of the Award for Excellence in Teaching at the California Institute of Technology for the years 1983–1984, 1992–1993, and 1993–1994. He also received the NSF's Presidential Young Investigator award in 1986. In 1989, he received the IEEE ASSP Senior Award for his paper on multirate perfect-reconstruction filter banks. In 1990 he was recipient of the S. K. Mitra Memorial Award from the Institute of Electronics and Telecommunications Engineers, India, for his joint paper in the IETE journal. In 2009 he was chosen to receive the IETE students journal award for his tutorial paper in the *IETE Journal of Education*. He was also the coauthor of a paper on linear-phase perfect reconstruction filter banks in the IEEE TRANSACTIONS ON SIGNAL PROCESSING, for which the first author (T. Nguyen) received the Young Outstanding Author award in 1993. He was elected Fellow of the IEEE in 1991. He received the 1995 F. E. Terman Award of the American Society for Engineering Education, sponsored by Hewlett Packard Co., for his contributions to engineering education, especially the book *Multirate Systems and Filter Banks* published by Prentice-Hall in 1993. He has given several plenary talks, including at the IEEE ISCAS-04, SAMPTA-01, EUSIPCO-98, SPCOM-95, and Asilomar-88 conferences on signal processing. He has been chosen a distinguished lecturer for the IEEE Signal Processing Society for the year 1996–1997. In 1999 he was chosen to receive the IEEE CAS Society's Golden Jubilee Medal. He is a recipient of the IEEE Signal Processing Society's Technical Achievement Award for the year 2002.